



TITLE:

# Estimation of a signal matrix for high-dimensional non-Gaussian data (Statistical Inference on Divergence Measures and Its Related Topics)

AUTHOR(S):

矢田, 和善; 青嶋, 誠

---

CITATION:

矢田, 和善 ...[et al]. Estimation of a signal matrix for high-dimensional non-Gaussian data (Statistical Inference on Divergence Measures and Its Related Topics). 数理解析研究所講究録 2016, 1999: 36-46; KJ00010266683.

ISSUE DATE:

2016-07

URL:

<http://hdl.handle.net/2433/224772>

RIGHT:

# Estimation of a signal matrix for high-dimensional non-Gaussian data

筑波大学・数理物質系 矢田 和善

Kazuyoshi Yata  
Institute of Mathematics  
University of Tsukuba

筑波大学・数理物質系 青嶋 誠

Makoto Aoshima  
Institute of Mathematics  
University of Tsukuba

*Abstract:* We consider the problem of recovering a low-rank signal matrix in high-dimensional situations. We first consider the conventional PCA to recover the signal matrix and show that the estimation of the signal matrix holds consistency properties under severe conditions. The conventional PCA is heavily subjected to the noise. In order to reduce the noise, we consider using the noise reduction (NR) methodology to recover the signal matrix and show that the estimation by the NR method improves the error rate of the conventional PCA effectively. We also apply the cross-data-matrix (CDM) methodology to recover the signal matrix and propose a new estimation of the signal matrix. We show that the proposed estimation by the CDM method performs well for high-dimensional non-Gaussian data.

*Key words and phrases:* Cross-data-matrix methodology, HDLSS, Large  $p$  small  $n$ , Noise-reduction methodology.

## 1 Introduction

In this paper, we address the problem of recovering an unknown  $d \times n$  low-rank matrix,  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n]$ .  $\mathbf{A}$  is called the signal matrix. Let  $r = \text{rank}(\mathbf{A})$ . We assume  $r$  ( $< \min\{d, n\}$ ) is fixed. Suppose we have a  $d \times n$  data matrix,  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]$ , where

$$\mathbf{X} = \sqrt{n}\mathbf{A} + \mathbf{W}. \quad (1)$$

Here,  $\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_n]$  is a  $d \times n$  noise matrix, where  $\mathbf{w}_j$ ,  $j = 1, \dots, n$ , are independent and identically distributed (i.i.d.) as a  $d$ -dimensional distribution with mean zero and covariance matrix  $\Sigma_W$  ( $\neq \mathbf{O}$ ). Note that  $\mathbf{x}_j = \sqrt{n}\mathbf{a}_j + \mathbf{w}_j$ ,  $j = 1, \dots, n$ ,

are i.i.d. Let  $\Sigma_A = \mathbf{A}\mathbf{A}^T$ . Then, it holds that  $E(\mathbf{X}\mathbf{X}^T)/n = \Sigma_A + \Sigma_W (= \Sigma, \text{ say})$ . Shabalin and Nobel [6] considered (1) in a high-dimensional setting, where the data dimension  $d$  and the sample size  $n$  increase at the same rate, i.e.  $n/d \rightarrow c > 0$ . They assumed that the elements of  $\mathbf{W}$  are i.i.d. normal random variables. We note that the conditions such as “ $n/d \rightarrow c > 0$ ” and the Gaussianity of the noise are often strict in real high-dimensional analyses. Yata and Aoshima [10] considered (1) in high-dimension, low-sample-size (HDLSS) settings without assuming those conditions.

The eigen-decomposition of  $\Sigma_W$  is given by  $\Sigma_W = \mathbf{U}_W \Lambda_W \mathbf{U}_W^T$ , where  $\Lambda_W$  is a diagonal matrix of eigenvalues,  $\lambda_{1(W)} \geq \dots \geq \lambda_{d(W)} (\geq 0)$ , and  $\mathbf{U}_W$  is an orthogonal matrix of the corresponding eigenvectors. Let  $\mathbf{W} = \mathbf{U}_W \Lambda_W^{1/2} \mathbf{Z}$ . Then,  $\mathbf{Z}$  is a  $d \times n$  sphered data matrix from a distribution with the identity covariance matrix. Here, we write  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_d]^T$  and  $\mathbf{z}_j = (z_{j1}, \dots, z_{jn})^T$ ,  $j = 1, \dots, d$ . Note that  $E(z_{jk}z_{j'k}) = 0$  ( $j \neq j'$ ) and  $\text{Var}(\mathbf{z}_j) = \mathbf{I}_n$ , where  $\mathbf{I}_n$  is the  $n$ -dimensional identity matrix. We assume that the fourth moments of each variable in  $\mathbf{Z}$  are uniformly bounded. The singular value decomposition of  $\mathbf{A}$  is given by

$$\mathbf{A} = \sum_{j=1}^r \lambda_{j(A)}^{1/2} \mathbf{u}_{j(A)} \mathbf{v}_{j(A)}^T,$$

where  $\lambda_{1(A)}^{1/2} \geq \dots \geq \lambda_{r(A)}^{1/2} (> 0)$  are singular values of  $\mathbf{A}$  and  $\mathbf{u}_{j(A)}$  (or  $\mathbf{v}_{j(A)}$ ) denotes a unit left- (or right-) singular vector corresponding to  $\lambda_{j(A)}^{1/2}$  ( $j = 1, \dots, r$ ). Note that  $\Sigma_A = \sum_{j=1}^r \lambda_{j(A)} \mathbf{u}_{j(A)} \mathbf{u}_{j(A)}^T$ . Also, note that  $\lambda_{j(A)}$ s depend not only on  $d$  but also on  $n$ . In this paper, we assume the following model.

$$\lim_{d \rightarrow \infty} \frac{\text{tr}(\Sigma_W^2)}{\lambda_{r(A)}^2} = 0 \quad \text{when } n \text{ is fixed or } n \rightarrow \infty. \quad (2)$$

The model (2) is a special case of the power spiked model given by Yata and Aoshima [9]. Murayama et al. [5] considered the estimation of  $\mathbf{A}$  for a special case of (2). When  $r \geq 2$ , we assume that  $\lambda_{j(A)}$ s are distinct in the sense that

$$\liminf_{d \rightarrow \infty} \frac{\lambda_{j(A)}}{\lambda_{j'(A)}} > 1 \quad \text{when } n \text{ is fixed or } n \rightarrow \infty \text{ for all } j < j' (\leq r).$$

The sample covariance matrix is given by  $\mathbf{S} = n^{-1} \mathbf{X}\mathbf{X}^T$ . We consider the dual sample covariance matrix defined by  $\mathbf{S}_D = n^{-1} \mathbf{X}^T \mathbf{X}$ . Let  $m = \min\{d, n\}$ . Note that  $\mathbf{S}_D$  and  $\mathbf{S}$  share non-zero eigenvalues and  $\text{rank}(\mathbf{S}) = \text{rank}(\mathbf{S}_D) \leq m$ . Let  $\hat{\lambda}_1 \geq \dots \geq \hat{\lambda}_m \geq 0$  be the eigenvalues of  $\mathbf{S}_D$ . The eigen-decompositions of  $\mathbf{S}$  and  $\mathbf{S}_D$  are given by  $\mathbf{S} = \sum_{j=1}^m \hat{\lambda}_j \hat{\mathbf{u}}_j \hat{\mathbf{u}}_j^T$  and  $\mathbf{S}_D = \sum_{j=1}^m \hat{\lambda}_j \hat{\mathbf{v}}_j \hat{\mathbf{v}}_j^T$ , where  $\hat{\mathbf{u}}_j$  (or  $\hat{\mathbf{v}}_j$ ) denotes a unit left- (or right-) singular vector of  $\mathbf{X}/n^{1/2}$  corresponding to  $\hat{\lambda}_j^{1/2}$ . Note that  $\hat{\mathbf{u}}_j$  can be calculated by  $\hat{\mathbf{u}}_j = (n\hat{\lambda}_j)^{-1/2} \mathbf{X} \hat{\mathbf{v}}_j$  from the fact that  $\mathbf{X}/n^{1/2} = \sum_{j=1}^m \hat{\lambda}_j^{1/2} \hat{\mathbf{u}}_j \hat{\mathbf{v}}_j^T$ .

Jung and Marron [4] investigated the inconsistency of the eigenvalues and eigenvectors of the sample covariance matrix for HDLSS data. Yata and Aoshima [8] gave consistent estimators for both the eigenvalues and eigenvectors together with the principal component (PC) scores by developing the *noise-reduction (NR) methodology*. Moreover, Yata and Aoshima [7] created a new principal component analysis (PCA) called the *cross-data-matrix (CDM) methodology* that is applicable to constructing an unbiased estimator in nonparametric settings. Aoshima and Yata [1] developed a variety of high-dimensional statistical inference by using the cross-data-matrix methodology. See Aoshima and Yata [2, 3] for a review covering this field of research.

In this paper, we consider the problem of recovering the signal matrix  $\mathbf{A}$  in high-dimensional settings. In Sections 2 and 3, we introduce the results of Yata and Aoshima [10]. In Section 2, we consider using the conventional PCA to recover  $\mathbf{A}$  and show that the estimation of  $\mathbf{A}$  holds consistency properties under severe conditions. In Section 3, we consider the noise reduction (NR) methodology by Yata and Aoshima [8] in (1) and apply it to recovering  $\mathbf{A}$ . We show that the estimation of  $\mathbf{A}$  by the NR method holds the consistency properties under mild conditions and improves the error rate of the conventional PCA. In Section 4, we consider the cross-data-matrix (CDM) methodology by Yata and Aoshima [7] in (1) and apply it to recovering  $\mathbf{A}$ . We show that the estimation of  $\mathbf{A}$  by the CDM method performs well for high-dimensional non-Gaussian data.

## 2 Estimation of the signal matrix by conventional PCA

In this section, we consider recovering the signal matrix  $\mathbf{A}$  by using the conventional PCA in high-dimensional settings such as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ . We reconstruct  $\mathbf{A}$  by using  $\hat{\lambda}_j$ s,  $\hat{\mathbf{u}}_j$ s and  $\hat{\mathbf{v}}_j$ s. We assume  $\hat{\mathbf{u}}_j^T \mathbf{u}_{j(A)} \geq 0$  and  $\hat{\mathbf{v}}_j^T \mathbf{v}_{j(A)} \geq 0$  for all  $j$  ( $\leq r$ ) without loss of generality.

We consider the following conditions when  $d \rightarrow \infty$  while  $n$  is fixed or  $n \rightarrow \infty$ :

$$(C-i) \quad \frac{\sum_{s,t=1}^d \lambda_{s(W)} \lambda_{t(W)} E\{(z_{sk}^2 - 1)(z_{tk}^2 - 1)\}}{n \lambda_{r(A)}^2} = o(1);$$

$$(C-ii) \quad \frac{\text{tr}(\Sigma_W)}{n \lambda_{r(A)}} = o(1).$$

**Remark 1.** We note that  $z_{1k}, \dots, z_{dk}$  ( $k = 1, \dots, n$ ) are independent when  $\mathbf{W}$  is Gaussian. Then, it holds that

$$\sum_{s,t=1}^d \lambda_{s(W)} \lambda_{t(W)} E\{(z_{sk}^2 - 1)(z_{tk}^2 - 1)\} = O\{\text{tr}(\Sigma_W^2)\}, \quad (3)$$

so that (C-i) holds under (2) when  $\mathbf{W}$  is Gaussian or  $z_{1k}, \dots, z_{dk}$  ( $k = 1, \dots, n$ ) are independent.

If (3) holds, (C-i) is met even when  $n$  is fixed. Let  $\kappa_j = \text{tr}(\Sigma_{\mathbf{W}})/(n\lambda_{j(A)})$  for  $j = 1, \dots, r$ . Yata and Aoshima [10] gave the following results.

**Theorem 2.1** ([10]). *Under (C-i), it holds that for  $j = 1, \dots, r$*

$$\begin{aligned} \frac{\hat{\lambda}_j}{\lambda_{j(A)}} &= 1 + \kappa_j + o_p(1), \quad \hat{\mathbf{u}}_j^T \mathbf{u}_{j(A)} = (1 + \kappa_j)^{-1/2} + o_p(1) \\ \text{and } \hat{\mathbf{v}}_j^T \mathbf{v}_{j(A)} &= 1 + o_p(1) \end{aligned}$$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

**Corollary 2.1** ([10]). *Under (C-i) and (C-ii), it holds that for  $j = 1, \dots, r$*

$$\frac{\hat{\lambda}_j}{\lambda_{j(A)}} = 1 + o_p(1) \quad \text{and} \quad \hat{\mathbf{u}}_j^T \mathbf{u}_{j(A)} = 1 + o_p(1)$$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

Note that  $\hat{\mathbf{v}}_j$ s hold the consistency property without (C-ii) contrary to  $\hat{\lambda}_j$ s and  $\hat{\mathbf{u}}_j$ s. Based on the theoretical background, we consider recovering the signal matrix  $\mathbf{A}$  by  $\hat{\mathbf{A}} = \sum_{i=1}^r \hat{\lambda}_i^{1/2} \hat{\mathbf{u}}_i \hat{\mathbf{v}}_i^T$ . Yata and Aoshima [10] discussed the choice of  $r$  in  $\hat{\mathbf{A}}$ . We define a loss function by

$$L(\hat{\mathbf{A}}|\mathbf{A}) = \|\hat{\mathbf{A}} - \mathbf{A}\|_F^2,$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. Let  $\psi = \text{tr}(\Sigma_{\mathbf{W}})/n$ . Then, Yata and Aoshima [10] gave the following results.

**Theorem 2.2** ([10]). *Under (C-i), it holds that*

$$L(\hat{\mathbf{A}}|\mathbf{A}) = r\psi + o_p(\lambda_{r(A)})$$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

**Remark 2.** *If (3) holds, Theorem 2.2 is claimed even when  $n$  is fixed.*

**Corollary 2.2** ([10]). *Under (C-i) and (C-ii), it holds that*

$$L(\hat{\mathbf{A}}|\mathbf{A}) = o_p(\lambda_{r(A)})$$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

From Theorem 2.2, if (C-ii) does not hold, the loss of  $\hat{\mathbf{A}}$  becomes  $r\text{tr}(\Sigma_{\mathbf{W}})/n$  asymptotically. In order to reduce the noise, we consider using the NR method to recover the signal matrix in Section 3.

### 3 Estimation of the signal matrix by NR method

In this section, we consider applying the *noise-reduction (NR) methodology* by Yata and Aoshima [8] to recover the signal matrix  $\mathbf{A}$ . By using the NR method, we obtain an estimator of  $\lambda_{j(A)}$  as

$$\hat{\lambda}_{j(r)} = \hat{\lambda}_j - \frac{\text{tr}(\mathbf{S}_D) - \sum_{i=1}^r \hat{\lambda}_i}{n - r} \quad (j = 1, \dots, r). \quad (4)$$

Note that the second term in (4) is an estimator of  $\psi$ . Then, Yata and Aoshima [10] gave the following result.

**Theorem 3.1** ([10]). *Under (C-i), it holds that for  $j = 1, \dots, r$*

$$\frac{\hat{\lambda}_{j(r)}}{\lambda_{j(A)}} = 1 + o_p(1)$$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

From Theorem 3.1,  $\hat{\lambda}_{j(r)}$  holds the consistency property without (C-ii). Remember that  $\hat{\lambda}_j$  requires (C-ii) to hold the consistency property.

We consider recovering  $\mathbf{A}$  by  $\hat{\mathbf{A}} = \sum_{i=1}^r \hat{\lambda}_{i(r)}^{1/2} \hat{\mathbf{u}}_i \hat{\mathbf{v}}_i^T$ . Yata and Aoshima [10] discussed the choice of  $r$  in  $\hat{\mathbf{A}}$ . Let

$$\delta_i = \mathbf{u}_{i(A)}^T \mathbf{W} \mathbf{v}_{i(A)} / (n \lambda_{i(A)})^{1/2} \quad \text{for } i = 1, \dots, r.$$

For the loss function by  $L(\hat{\mathbf{A}}|\mathbf{A}) = \|\hat{\mathbf{A}} - \mathbf{A}\|_F^2$ , Yata and Aoshima [10] gave the following results.

**Theorem 3.2.** *Under (C-i), it holds that*

$$L(\hat{\mathbf{A}}|\mathbf{A}) = 2 \sum_{i=1}^r \lambda_{i(A)} (1 + \delta_i) \left( 1 - \frac{1 + \delta_i}{(1 + \kappa_i + 2\delta_i)^{1/2}} \right) + o_p(\lambda_{r(A)})$$

and  $\delta_i = o_p\{(\lambda_{r(A)}/\lambda_{i(A)})^{1/2}\}$  for  $i = 1, \dots, r$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

**Remark 3.** *If (3) holds, Theorem 3.2 is claimed even when  $n$  is fixed.*

**Corollary 3.1** ([10]). *Under (C-i) and (C-ii), it holds that*

$$L(\hat{\mathbf{A}}|\mathbf{A}) = o_p(\lambda_{r(A)})$$

as  $d \rightarrow \infty$  either when  $n$  is fixed or  $n \rightarrow \infty$ .

From Theorems 2.2 and 3.2, we compare  $2\lambda_{i(A)}\{1 - 1/(1 + \kappa_i)^{1/2}\}$  with  $\psi$  ( $= \lambda_{i(A)}\kappa_i$ ) by noting  $\delta_i = o_p(1)$ . It holds that  $2\{1 - 1/(1 + \kappa_i)^{1/2}\} < \kappa_i$  ( $i = 1, \dots, r$ ) for any  $\kappa_i > 0$ , so that  $L(\hat{\mathbf{A}}|\mathbf{A})$  is smaller than  $L(\hat{\mathbf{A}}|\mathbf{A})$  asymptotically. Thus,  $\hat{\mathbf{A}}$  improves the loss of  $\hat{\mathbf{A}}$ .

## 4 Estimation of the signal matrix by CDM method

In this section, we consider applying the *cross-data-matrix (CDM) methodology* by Yata and Aoshima [7] to recover the signal matrix  $\mathbf{A}$ .

Let  $n_1 = \lceil n/2 \rceil$  and  $n_2 = n - n_1$ , where  $\lceil x \rceil$  denotes the smallest integer  $\geq x$ . We assume  $n_2 \geq r$ . Let  $\mathbf{X}_1 = [\mathbf{x}_1, \dots, \mathbf{x}_{n_1}]$  and  $\mathbf{X}_2 = [\mathbf{x}_{n_1+1}, \dots, \mathbf{x}_n]$ . We define

$$\mathbf{S}_{D(1)} = (n_1 n_2)^{-1/2} \mathbf{X}_1^T \mathbf{X}_2.$$

Let  $m_0 = \min\{d, n_2\}$ . The singular value decomposition of  $\mathbf{S}_{D(1)}$  is given by

$$\mathbf{S}_{D(1)} = \sum_{j=1}^{m_0} \tilde{\lambda}_j \tilde{\mathbf{v}}_{1j} \tilde{\mathbf{v}}_{2j}^T,$$

where  $\tilde{\lambda}_j \geq \dots \geq \tilde{\lambda}_{m_0} (\geq 0)$  denote singular values of  $\mathbf{S}_{D(1)}$  and  $\tilde{\mathbf{v}}_{1j}$  (or  $\tilde{\mathbf{v}}_{2j}$ ) denotes a unit left- (or right-) singular vector corresponding to  $\tilde{\lambda}_j$  ( $j = 1, \dots, m_0$ ). Let  $\mathbf{A}_1 = [\mathbf{a}_1, \dots, \mathbf{a}_{n_1}]$ ,  $\mathbf{A}_2 = [\mathbf{a}_{n_1+1}, \dots, \mathbf{a}_n]$ ,  $\mathbf{W}_1 = [\mathbf{w}_1, \dots, \mathbf{w}_{n_1}]$  and  $\mathbf{W}_2 = [\mathbf{w}_{n_1+1}, \dots, \mathbf{w}_n]$ . Then, we write that

$$\mathbf{X}_i = \sqrt{n} \mathbf{A}_i + \mathbf{W}_i, \quad i = 1, 2.$$

Let  $\mathbf{v}_{j(A)} = (\mathbf{v}_{1j(A)}^T, \mathbf{v}_{2j(A)}^T)^T$  for  $j = 1, \dots, r$ , where  $\mathbf{v}_{ij(A)} \in \mathbb{R}^{n_i}$ . Note that  $\|\mathbf{v}_{1j(A)}\|^2 + \|\mathbf{v}_{2j(A)}\|^2 = \|\mathbf{v}_{j(A)}\|^2 = 1$  for  $j = 1, \dots, r$ , where  $\|\cdot\|$  denotes the Euclidean norm. Then, we write that

$$\mathbf{A}_i = \sum_{j=1}^r \lambda_{j(A)}^{1/2} \mathbf{u}_{j(A)} \mathbf{v}_{ij(A)}^T, \quad i = 1, 2.$$

Hereafter, we assume that

$$\limsup_{m \rightarrow \infty} \frac{\lambda_{1(A)}}{\lambda_{r(A)}} < \infty \quad \text{and} \quad \|\mathbf{v}_{ij(A)}\|^2 = 1/2 + o(1) \quad \text{as } m \rightarrow \infty \text{ for all } i, j. \quad (5)$$

We assume  $\tilde{\mathbf{v}}_{ij}^T \mathbf{v}_{ij(A)} \geq 0$  for all  $i, j$  without loss of generality. Then, we have the following result.

**Theorem 4.1.** *It holds that for  $j = 1, \dots, r$*

$$\frac{\tilde{\lambda}_j}{\lambda_{j(A)}} = 1 + o_p(1) \quad \text{and} \quad \sqrt{2} \tilde{\mathbf{v}}_{ij}^T \mathbf{v}_{ij(A)} = 1 + o_p(1), \quad i = 1, 2,$$

as  $m \rightarrow \infty$ .

From Theorem 4.1,  $\tilde{\lambda}_j$  holds the consistency property without (C-i) and (C-ii). Remember that  $\hat{\lambda}_{j(r)}$  requires (C-i) to hold the consistency property.

Next, we consider estimation of  $\mathbf{u}_{j(A)}$ s by using the CDM methodology. Let  $\tilde{\mathbf{u}}_{ij} = (n_i \tilde{\lambda}_j)^{-1/2} \mathbf{X}_i \tilde{\mathbf{v}}_{ij}$ ,  $i = 1, 2$ . Then, we have the following result.

**Theorem 4.2.** *It holds that for  $j = 1, \dots, r$*

$$\tilde{\mathbf{u}}_{ij}^T \mathbf{u}_{j(A)} = 1 + o_p(1), \quad i = 1, 2,$$

as  $m \rightarrow \infty$ .

Let  $\tilde{\mathbf{u}}_{ij(*)} = \tilde{\mathbf{u}}_{ij}/\|\tilde{\mathbf{u}}_{ij}\|$  for all  $i, j$ . Let  $\tilde{\mathbf{A}}_i = \sum_{j=1}^r \tilde{\lambda}_j^{1/2} \tilde{\mathbf{u}}_{ij(*)} \tilde{\mathbf{v}}_{ij}^T / \sqrt{2}$  for  $i = 1, 2$ . We consider recovering  $\mathbf{A}$  by  $\tilde{\mathbf{A}} = [\tilde{\mathbf{A}}_1, \tilde{\mathbf{A}}_2]$ . From Theorems 4.1 and 4.2, we expect that  $\tilde{\mathbf{A}}$  performs well for high-dimensional non-Gaussian data. In Section 5, we examine the performance of  $\tilde{\mathbf{A}}$  with the help of numerical simulations.

## 5 Simulations

In this section, we give numerical comparisons of  $\hat{\mathbf{A}}$ ,  $\dot{\mathbf{A}}$  and  $\tilde{\mathbf{A}}$ .

We set  $d = 2^t$ ,  $t = 4, \dots, 10$ ,  $n = 15$ ,  $r = 3$ ,  $\Sigma_A = \text{diag}(\lambda_{1(A)}, \lambda_{2(A)}, \lambda_{3(A)}, 0, \dots, 0)$  with  $(\lambda_{1(A)}, \lambda_{2(A)}, \lambda_{3(A)}) = (d/5, d/15, d/45)$  and  $\Sigma_W = (0.3^{|i-j|^{1/3}})$ . Note that  $\text{tr}(\Sigma_W) = d$ . We considered three cases:

- (a)  $\mathbf{w}_k$ s are i.i.d. as a  $d$ -variate normal distribution,  $N_d(\mathbf{0}, \Sigma_W)$  with mean zero and covariance matrix  $\Sigma_W$ ;
- (b)  $\mathbf{w}_k$ s are i.i.d. as a  $d$ -variate  $t$ -distribution,  $t_d(\mathbf{0}, \Sigma_W, \nu)$  with mean zero, covariance matrix  $\Sigma_W$  and degrees of freedom  $\nu = 10$ ;
- (c)  $\mathbf{w}_k$ s are i.i.d. as a  $d$ -variate  $t$ -distribution,  $t_d(\mathbf{0}, \Sigma_W, \nu)$  with  $\nu = 30$ .

Let  $F(\mathbf{M}) = \|\mathbf{M} - \mathbf{A}\|_F^2/d$  for any  $d \times n$  matrix,  $\mathbf{M}$ . The findings were obtained by averaging the outcomes from 2000 ( $= K$ , say) replications. Under a fixed scenario, suppose that the  $k$ -th replication ends with estimates,  $F(\hat{\mathbf{A}})_k$ ,  $F(\dot{\mathbf{A}})_k$  and  $F(\tilde{\mathbf{A}})_k$ , for  $k = 1, \dots, K$ . Let us simply write  $\hat{F} = K^{-1} \sum_{k=1}^K F(\hat{\mathbf{A}})_k$ ,  $\dot{F} = K^{-1} \sum_{k=1}^K F(\dot{\mathbf{A}})_k$  and  $\tilde{F} = K^{-1} \sum_{k=1}^K F(\tilde{\mathbf{A}})_k$ . We also considered the Monte Carlo variability. Let  $\text{var}(\hat{F}) = (K-1)^{-1} \sum_{k=1}^K (F(\hat{\mathbf{A}})_k - \hat{F})^2$ ,  $\text{var}(\dot{F}) = (K-1)^{-1} \sum_{k=1}^K (F(\dot{\mathbf{A}})_k - \dot{F})^2$  and  $\text{var}(\tilde{F}) = (K-1)^{-1} \sum_{k=1}^K (F(\tilde{\mathbf{A}})_k - \tilde{F})^2$ . Figure 1 shows the behaviors of  $(\hat{F}, \dot{F}, \tilde{F})$  and  $(\text{var}(\hat{F}), \text{var}(\dot{F}), \text{var}(\tilde{F}))$  for (a), (b) and (c).

We observed that the NR method and the CDM method give more preferable performances compared to the conventional PCA. It seems that the NR method performs better than the CDM method for (a). Note that  $t_d(\mathbf{0}, \Sigma_W, \nu) \Rightarrow N_d(\mathbf{0}, \Sigma_W)$  as  $\nu \rightarrow \infty$ . When  $\nu = 10$ , the NR method seems not to give a feasible estimation. This is probably because  $\nu = 10$  is not large enough for  $\mathbf{W}$  to satisfy (C-i). On the other hand, the CDM method does not require (C-i). As observed in Figure 2, the CDM method seems to perform well even when  $\nu = 10$ .

## A Appendix

Throughout, let  $\mathbf{e}_{in_i} = (e_{i1}, \dots, e_{in_i})^T$ ,  $i = 1, 2$ , be arbitrary unit random vectors. Let  $\mathbf{z}_{1j} = (z_{j1}, \dots, z_{jn_1})^T$  and  $\mathbf{z}_{2j} = (z_{jn_1+1}, \dots, z_{jn})^T$ ,  $j = 1, \dots, d$ .



**Lemma 1.** *It holds that as  $m \rightarrow \infty$*

$$\mathbf{e}_{1n_1}^T \frac{\mathbf{W}_1^T \mathbf{W}_2}{(n_1 n_2)^{1/2} \lambda_{r(A)}} \mathbf{e}_{2n_2} = o_p(1)$$

under (2).

*Proof.* We write that

$$\mathbf{e}_{1n_1}^T \frac{\mathbf{W}_1^T \mathbf{W}_2}{(n_1 n_2)^{1/2} \lambda_{r(A)}} \mathbf{e}_{2n_2} = \mathbf{e}_{1n_1}^T \frac{\sum_{j=1}^d \lambda_j(W) \mathbf{z}_{1j} \mathbf{z}_{2j}^T}{(n_1 n_2)^{1/2} \lambda_{r(A)}} \mathbf{e}_{2n_2}.$$

Then, by using Lemma 4 given in Yata and Aoshima [9], we can conclude the result.  $\square$

**Lemma 2.** *It holds that as  $m \rightarrow \infty$*

$$\mathbf{e}_{1n_1}^T \frac{\mathbf{A}_1^T \mathbf{W}_2}{n_2^{1/2} \lambda_{r(A)}} \mathbf{e}_{2n_2} = o_p(1) \quad \text{and} \quad \mathbf{e}_{1n_1}^T \frac{\mathbf{W}_1^T \mathbf{A}_2}{n_1^{1/2} \lambda_{r(A)}} \mathbf{e}_{2n_2} = o_p(1)$$

under (2) and (5).

*Proof.* We first consider the first result of Lemma 2. We note that

$$|\mathbf{e}_{1n_1}^T \mathbf{A}_1^T \mathbf{W}_2 \mathbf{e}_{2n_2}| \leq \sum_{i=1}^r \lambda_{i(A)}^{1/2} |\mathbf{u}_{i(A)}^T \mathbf{W}_2 \mathbf{e}_{2n_2}|. \quad (6)$$

We write that  $\mathbf{u}_{i(A)}^T \mathbf{W}_2 \mathbf{e}_{2n_2} = \sum_{k=1}^{n_2} e_{2k} \mathbf{w}_{n_1+k}^T \mathbf{u}_{i(A)}$ . Note that  $\lambda_1(W) = o(\lambda_{r(A)})$  as  $m \rightarrow \infty$  from (2). By using Markov's inequality, for any  $\tau > 0$  and  $i = 1, \dots, r$ , we have that as  $m \rightarrow \infty$

$$\begin{aligned} P\left(\sum_{k=1}^{n_2} (\mathbf{w}_{n_1+k}^T \mathbf{u}_{i(A)})^2 / n_2 \geq \tau \lambda_{r(A)}\right) &\leq \frac{E\{\sum_{k=1}^{n_2} (\mathbf{w}_{n_1+k}^T \mathbf{u}_{i(A)})^2\}}{\tau n_2 \lambda_{r(A)}} = \frac{\mathbf{u}_{i(A)}^T \Sigma_W \mathbf{u}_{i(A)}}{\tau \lambda_{r(A)}} \\ &\leq \frac{\lambda_1(W)}{\tau \lambda_{r(A)}} = o(1) \end{aligned}$$

from the fact that  $\mathbf{u}_{i(A)}^T \Sigma_W \mathbf{u}_{i(A)} \leq \lambda_1(W)$ . Then, by noting that

$$\begin{aligned} \left| \sum_{k=1}^{n_2} e_{2k} (\mathbf{w}_{n_1+k}^T \mathbf{u}_{i(A)}) / n_2^{1/2} \right| &\leq \left\{ \sum_{k=1}^{n_2} e_{2k}^2 \right\}^{1/2} \left\{ \sum_{k=1}^{n_2} (\mathbf{w}_{n_1+k}^T \mathbf{u}_{i(A)})^2 / n_2 \right\}^{1/2} \\ &= \left\{ \sum_{k=1}^{n_2} (\mathbf{w}_{n_1+k}^T \mathbf{u}_{i(A)})^2 / n_2 \right\}^{1/2} = o_p(\lambda_{r(A)}^{1/2}), \end{aligned} \quad (7)$$

from (6), we can conclude the first result. Similarly, we can conclude the second result. The proof is completed.  $\square$

*Proof of Theorem 4.1.* We write that for  $j = 1, \dots, r$

$$\frac{\tilde{\lambda}_j}{\lambda_{j(A)}} = \tilde{\mathbf{v}}_{1j}^T \frac{\mathbf{S}_{D(1)}}{\lambda_{j(A)}} \tilde{\mathbf{v}}_{2j} = \tilde{\mathbf{v}}_{1j}^T \frac{(n^{1/2} \mathbf{A}_1 + \mathbf{W}_1)^T (n^{1/2} \mathbf{A}_2 + \mathbf{W}_2)}{(n_1 n_2)^{1/2} \lambda_{j(A)}} \tilde{\mathbf{v}}_{2j}. \quad (8)$$

Then, by combining (8) with Lemmas 1 and 2, under (2) and (5), it holds that for  $j = 1, \dots, r$

$$\begin{aligned} \frac{\tilde{\lambda}_j}{\lambda_{j(A)}} &= 2\{1 + o(1)\} \tilde{\mathbf{v}}_{1j}^T \frac{\mathbf{A}_1^T \mathbf{A}_2}{\lambda_{j(A)}} \tilde{\mathbf{v}}_{2j} + o_p(1) \\ &= 2\{1 + o(1)\} \tilde{\mathbf{v}}_{1j}^T \frac{\sum_{j=1}^r \lambda_{j(A)} \mathbf{v}_{1j(A)} \mathbf{v}_{2j(A)}^T}{\lambda_{j(A)}} \tilde{\mathbf{v}}_{2j} = 1 + o_p(1). \end{aligned}$$

as  $m \rightarrow \infty$ . Thus, we have that

$$2^{1/2} \tilde{\mathbf{v}}_{ij}^T \mathbf{v}_{ij(A)} = 1 + o_p(1) \quad \text{for all } i, j. \quad (9)$$

It concludes the results.  $\square$

*Proof of Theorem 4.2.* We write that for all  $i, j$

$$\mathbf{u}_{j(A)}^T \tilde{\mathbf{u}}_{ij} = \frac{n^{1/2} \lambda_{j(A)}^{1/2} \mathbf{v}_{ij(A)}^T + \mathbf{u}_{j(A)}^T \mathbf{W}_i}{(n_i \tilde{\lambda}_j)^{1/2}} \tilde{\mathbf{v}}_{ij}.$$

Then, from Theorem 4.1, (7) and (9) under (2) and (5), it holds that for all  $i, j$

$$\mathbf{u}_{j(A)}^T \tilde{\mathbf{u}}_{ij} = 1 + o_p(1)$$

as  $m \rightarrow \infty$ . It concludes the result.  $\square$

## Acknowledgements

Research of the first author was partially supported by Grant-in-Aid for Young Scientists (B), Japan Society for the Promotion of Science (JSPS), under Contract Number 26800078. Research of the second author was partially supported by Grants-in-Aid for Scientific Research (A) and Challenging Exploratory Research, JSPS, under Contract Numbers 15H01678 and 26540010.

## References

- [1] M. Aoshima, K. Yata, Two-stage procedures for high-dimensional data, Sequential Anal. (*Editor's special invited paper*) 30 (2011) 356-399.

- [2] M. Aoshima, K. Yata, Invited review article: Statistical inference in high-dimension, low-sample-size settings, *Sugaku* 65 (2013) 225-247.
- [3] M. Aoshima, K. Yata, The JSS research prize lecture: Effective methodologies for high-dimensional data, *J. Japan Statist. Soc. Ser. J* 43 (2013) 123-150.
- [4] S. Jung, J. S. Marron, PCA consistency in high dimension, low sample size context, *Ann. Statist.* 37 (2009) 4104-4130.
- [5] W. Murayama, K. Yata, M. Aoshima, Reconstruction of a signal matrix for high-dimension, low-sample-size data, *RIMS Koukyuroku* 1954 (2015) 23-31.
- [6] A. Shabalin, A. Nobel, Reconstruction of a low-rank matrix in the presence of Gaussian noise, *J. Multivariate Anal.* 118 (2013) 67-76.
- [7] K. Yata, M. Aoshima, Effective PCA for high-dimension, low-sample-size data with singular value decomposition of cross data matrix, *J. Multivariate Anal.* 101 (2010) 2060-2077.
- [8] K. Yata, M. Aoshima, Effective PCA for high-dimension, low-sample-size data with noise reduction via geometric representations, *J. Multivariate Anal.* 105 (2012) 193-215.
- [9] K. Yata, M. Aoshima, PCA consistency for the power spiked model in high-dimensional settings, *J. Multivariate Anal.* 122 (2013) 334-354.
- [10] K. Yata, M. Aoshima, Reconstruction of a high-dimensional low-rank matrix, *Electron. J. Stat.*, accepted.

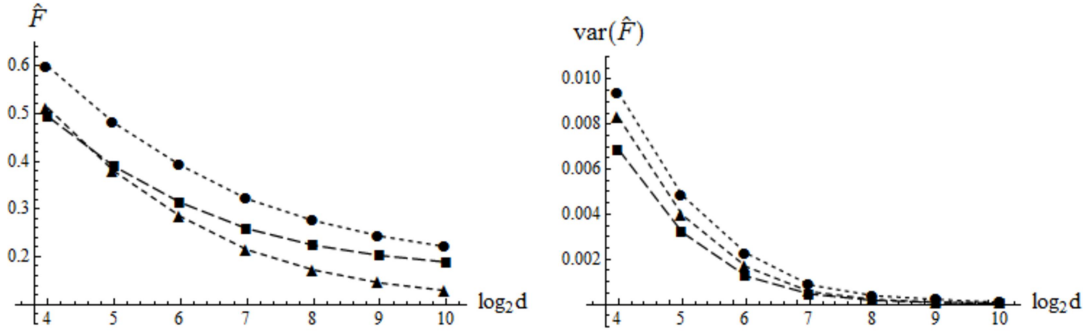
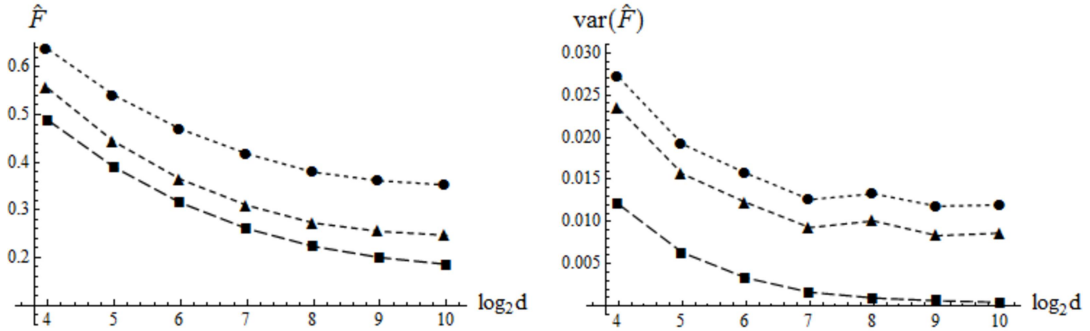
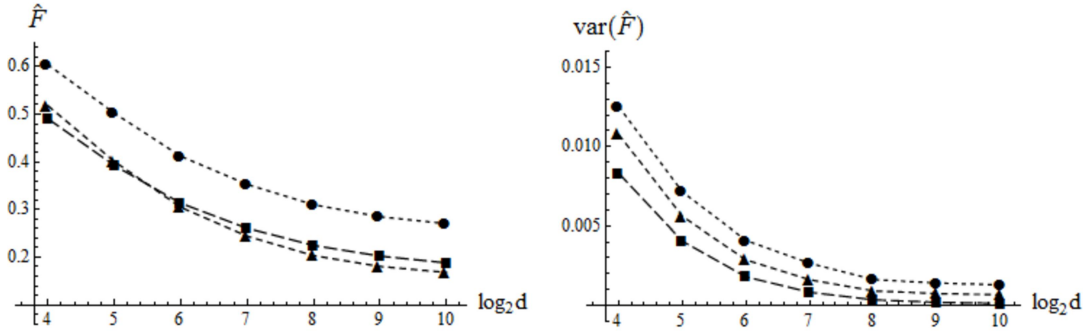
(a)  $N_d(\mathbf{0}, \Sigma_W)$ (b)  $t_d(\mathbf{0}, \Sigma_W, \nu)$  with  $\nu = 10$ (c)  $t_d(\mathbf{0}, \Sigma_W, \nu)$  with  $\nu = 30$ 

Figure 1: The behaviors of three estimates,  $\hat{A}$  denoted by  $\bullet$ ,  $\tilde{A}$  denoted by  $\blacktriangle$  and  $\tilde{\tilde{A}}$  denoted by  $\blacksquare$ . The values of  $\hat{F}$ ,  $\tilde{F}$  and  $\tilde{\tilde{F}}$  are given in the left panels and their sample variances,  $\text{var}(\hat{F})$ ,  $\text{var}(\tilde{F})$  and  $\text{var}(\tilde{\tilde{F}})$ , are given in the right panels.